

Language Models are Symbolic Learners in Arithmetic

Chunyuan Deng¹ Zhiqi Li² Roy Xie³ Ruidi Chang¹ Hanjie Chen¹

¹Rice University ²Georgia Tech ³Duke University

{chunyuan.deng,hanjie}@rice.edu

Abstract

Large Language Models (LLMs) are thought to struggle with arithmetic learning due to the inherent differences between language modeling and numerical computation, but concrete evidence has been lacking. This work responds to this claim through a two-side experiment. We first investigate whether LLMs leverage partial products during arithmetic learning. We find that although LLMs can identify some partial products after learning, they fail to leverage them for arithmetic tasks, conversely. We then explore how LLMs approach arithmetic symbolically by breaking tasks into subgroups, hypothesizing that difficulties arise from subgroup complexity and selection. Our results show that when subgroup complexity is fixed, LLMs treat a collection of different arithmetic operations similarly. By analyzing position-level accuracy across different training sizes, we further observe that it follows a U-shaped pattern: LLMs quickly learn the easiest patterns at the first and last positions, while progressively learning the more difficult patterns in the middle positions. This suggests that LLMs select subgroup following an easy-to-hard paradigm during learning. Our work confirms that LLMs are pure symbolic learners in arithmetic tasks and underscores the importance of understanding them deeply through subgroup-level quantification.

1 Introduction

Modern math benchmarks like MATH (Hendrycks et al., 2021) and GSM8K (Cobbe et al., 2021) have been rapidly saturated due to the advancements of frontier language models like GPT-4o (OpenAI et al., 2024) and Claude (Anthropic, 2024). This trend is driving the assessment of these models toward more challenging tasks, such as Olympic Math. However, it has been observed that even the most advanced language models struggle with basic arithmetic, such as 5-digit multiplication (Yang

et al., 2023). This notable gap raises questions about the underlying mechanisms behind arithmetic in language models. And some hypothesize (Boguraev et al., 2024) that this difficulty stems from the fact that mathematical calculation differs fundamentally from autoregressive language modeling.

Previous research on this topic has primarily focused on causal abstraction to identify model components such as key attention layers (Stolfo et al., 2023) or attention heads (Zhang et al., 2024) responsible for arithmetic learning. While these studies provide valuable insights into the components involved in mathematical operations, they fall short of explaining why frontier models continue to struggle with certain arithmetic tasks. For instance, causal attribution reveals that the 14th attention head in the 19th layer is responsible for the operation $37 + 14 = 51$. However, it cannot explain why the model handles $37 + 14$ successfully but fails on 37×14 . This observation suggests that there is still room to explore arithmetic learning from alternative perspectives.

To achieve this, we approach this problem from two sides (shown in Figure 1). First, we examine whether LLMs leverage partial products in arithmetic learning tasks. Then we explore whether, and more importantly, *how* LLMs handle arithmetic in a purely symbolic manner. Specifically, we decompose the task into subgroup level, hypothesizing the task can be understood through two aspects: subgroup complexity and subgroup selection.

For partial products, we considered four distinct methods for multiplication calculation: standard multiplication, repetitive addition, the lattice method, and Egyptian multiplication (detailed in §5.1), each generating distinct paths for partial product computation. We first fine-tune LLMs on multiplication tasks. After fine-tuning, models improved in identifying nearly all partial products, but explicit training on partial products did not

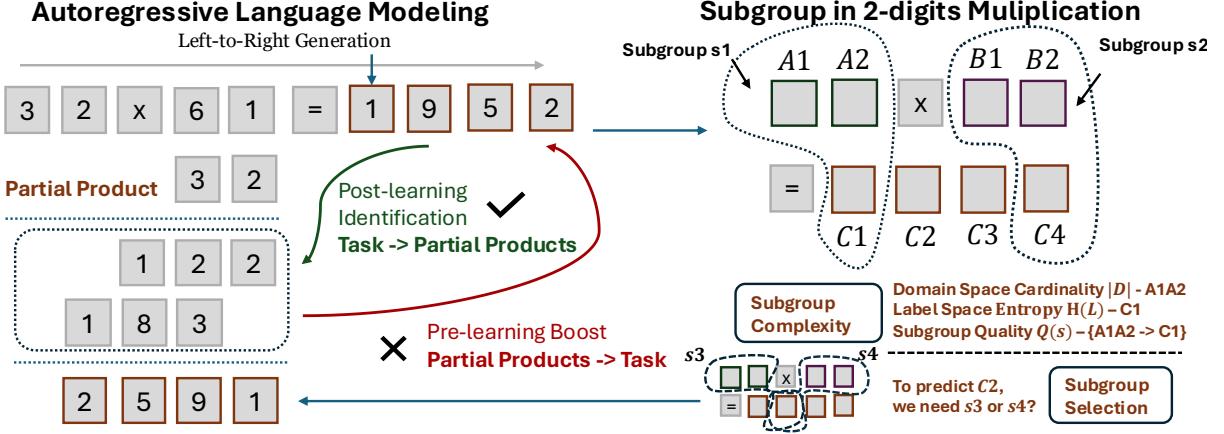


Figure 1: Fundamental structure of the paper. We begin by investigating partial products and proceed to a detailed examination at the subgroup level to understand the mechanism in a symbolic manner.

enhance their ability to perform multiplications, instead. These findings suggest that large language models may not leverage partial products to perform arithmetic calculations. The increased ability to identify partial products seems to be a by-product of subgroup-level symbol learning rather than a mechanism the models inherently use for computation.

We then delve into a more fine-grained level by examining subgroups to understand their basic complexity. Subgroup complexity refers to the intrinsic complexity of the subgroup itself. We propose that this is related to the domain space cardinality $|D|$, label space entropy $H(L)$, and subgroup quality $Q(s)$. Here, $|D|$ represents the maximum size of training data, $H(L)$ represents the variability in the label space, and $Q(s)$ refers to how deterministically one subgroup can map from the domain to the label space. While $|D|$, which strongly correlated with training size, is self-evident to influence learning, we will focus on label space entropy $H(L)$ first then discuss $Q(s)$ in the subgroup selection section. We observe influence of $H(L)$ by introducing perturbations to addition and multiplication by applying translation (Δc) and scaling ($\times \lambda$), while maintaining the total entropy the same across all output digits. Our findings show that LLMs have nearly the same accuracy across different perturbation magnitudes. Furthermore, when we intentionally reduce the entropy of the label space via modular operations, we observe an increase in accuracy as entropy decreases. These experiments confirm that label space entropy is an effective measure for quantifying task complexity and validates the hypothesis that LLMs operate as symbol-level pattern finders.

Subgroup selection refers to the process by

which LLMs identify the correct subgroup during training, specifically by selecting the appropriate mapping between a token subset in the input domain and a corresponding token subset in the output label space. To investigate this, we curated four distinct training sets to analyze position-level accuracy throughout the learning process. Consistent patterns emerged: position-level accuracy exhibited a U-shaped curve, achieving near-perfect accuracy (close to 100%) at both the beginning and end of digits with high $Q(s)$, but dropping significantly (to below 10%) across the intermediate digits with low $Q(s)$. These results suggest that LLMs apply a selection mechanism driven by high-to-low subgroup quality, providing further evidence of their symbolic learning behavior.

Our work confirms that LLMs do not truly perform calculations in arithmetic learning. Instead, they approach it in a purely symbolic manner. We then provide a systematic framework to examine this mechanism from subgroup complexity and subgroup selection. Our research emphasizes the pivotal role of label space entropy in the convergence stage and the impact of subgroup quality on learning dynamics, highlighting the importance of deeply understanding arithmetic through subgroup-level quantification symbolically.

2 Related Work

Mathematical Learning in LLMs Mathematical reasoning has been a longstanding area of research in natural language processing (NLP) (Kushman et al., 2014; Huang et al., 2016; Wang et al., 2017; Thawani et al., 2021; Sundaram et al., 2022; Guo et al., 2024). With recent advances in LLM, numerous studies have sought to improve their mathematical reasoning abilities through various

techniques, including data annealing (Dubey et al., 2024), continued pretraining (Lewkowycz et al., 2022), fine-tuning (Yue et al., 2023; Liu et al., 2023), prompting (Wei et al., 2023; Wang et al., 2023), and inference-time computation (Zhou et al., 2023a; Wu et al., 2024a). However, LLMs still face challenges with basic calculations and remain vulnerable to adversarial examples or perturbations, where minor changes in problems can result in incorrect answers (Zhou et al., 2023b; Xie et al., 2024). Most research on LLMs’ mathematical reasoning focuses on math word problems, where natural language is involved in the questions (Hendrycks et al., 2021; Cobbe et al., 2021; Arora et al., 2023; Zhao et al., 2024a,b).

Arithmetic Learning in Transformer Several previous efforts have aimed to improve arithmetic learning in LLMs. Lee et al. (2023) trained a 10.6M NanoGPT (Karpathy, 2022) model to learn arithmetic by carefully curating the data format, explicitly expanding each step using a method termed Scratchpad, which achieved remarkable performance compared to GPT-2 XL (Radford et al., 2019). Yang et al. (2023) fine-tuned MathGLM with a sufficient training dataset, demonstrating its capability to solve 5-digit multiplication. Deng et al. (2023, 2024) further advanced this field by internalizing the CoT process, hiding detailed steps in a scheduled manner, enabling GPT-2 small to solve 9-digit multiplication after multiple training runs.

Research on understanding arithmetic primarily stems from the interpretability community. The core idea is to identify *causal correlations* between model components and outputs. Stolfo et al. (2023) identify key attention layers responsible for arithmetic learning using causal mediation analysis (CMA), a weight perturbation method that observes changes in output. Similarly, Hanna et al. (2023) and Wu et al. (2024b) explore causal abstraction concepts at different model scales, specifically 0.1B and 7B parameters, respectively. More recently, Zhang et al. (2024) employed an attention attribution to isolate a small subset of attention heads and fine-tune them for improved performance at a lower cost. While these studies have made progress in understanding how LMs perform arithmetic at a component level, there remains a gap in understanding the learning mechanisms from a purely symbolic perspective. Our research aims to contribute to this missing gap in a systematic manner.

3 Preliminaries

In this section, we present the preliminaries of basic autoregressive language modeling, along with algebraic structure and arithmetic learning.

Autoregressive Language Modeling An autoregressive (AR) language model predicts the next token in a sequence based on the previously observed tokens. Formally, given a sequence of tokens $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$, the probability of the sequence is decomposed using the chain rule of probability as:

$$P(\mathbf{x}) = \prod_{t=1}^T P(x_t|x_1, x_2, \dots, x_{t-1}), \quad (1)$$

where $P(x_t|x_1, x_2, \dots, x_{t-1})$ represents the conditional probability of token x_t given all previous tokens. In autoregressive modeling, the next token is sampled from the conditional distribution $P(x_t|x_1, x_2, \dots, x_{t-1})$ until the end of the sequence is reached.

Algebraic Structure In our setting, we employ the algebraic structure known as a *ring*, which provides a formal framework for the arithmetic operations of *addition* and *multiplication*. A ring $(R, +, \cdot)$ is defined by:

- A set R (domain) of elements. Specifically, the domain R in our task is the set of integers \mathbb{Z} .
- Two binary operations, addition and multiplication $f(a, b) = c : R \times R \rightarrow R$. Specifically, we use $A_1A_2 + B_1B_2 = C_1C_2$ to represent 2-digit addition. We use $A_1A_2 \times B_1B_2 = C_1C_2C_3C_4$ to represent 2-digit multiplication.

In our case, for all $a, b \in R$, there exists a unique element $a + b \in R$, representing the sum of a and b . Similarly, there exists a unique element $a \cdot b \in R$, representing the product of a and b .

Arithmetic Learning Let \mathcal{M} denote a pretrained autoregressive language model with weights \mathbf{w} . We define an arithmetic task \mathcal{T} as a function learning problem where the goal is to predict numerical outputs based on arithmetic expressions. The training dataset for this task is given by:

$$\mathbb{D}_{\text{train}} = \{(a^{(k)}, b^{(k)}, c^{(k)})\}_{k=1}^N$$

where $c^{(k)} = f(a^{(k)}, b^{(k)})$ and $f(\cdot)$ represents a binary operator, N denotes the number of data points. In this context, $a^{(k)}$ and $b^{(k)}$ are input operands, and $c^{(k)}$ is the corresponding output, which is computed using the operator f (e.g., addition, multiplication, etc.)

4 Experiment Setup

In this section, we will detail our experiment setup. Unless otherwise specified, the same setup will be used in the following section.

Domain We select addition and multiplication as the fundamental operations for our experiments following previous work (Lee et al., 2023; Deng et al., 2023, 2024).

Model To investigate arithmetic learning at scale, we selected two open-source LLMs, Gemma-2-2B (Team et al., 2024) and Llama-3.1-8B (Dubey et al., 2024). Both models are top performers in their respective categories and excel at language-related tasks. We did not choose GPT-4o (OpenAI et al., 2024) or other proprietary LLMs due to concerns that they may internally integrate function calling (e.g., invoking APIs or executing Python programs), which could affect the experimental setup. Training details is included at Appendix A.1.

Conventional Data Format We directly train the model to predict the output (e.g., 130) given the input operands and the operator (e.g., 13×10). We add one space between each digit to ensure tokens are split into individual digits. We do not consider chain-of-thought (CoT) (Wei et al., 2023) or other prompting strategies to enforce the model to focus on arithmetic learning. We include ablations with respect to data format in Appendix A.2.

5 Are Large Language Models Implicit Calculators?

In this section, we explore whether LLMs utilize partial products to enhance their arithmetic calculation capabilities, particularly in the context of *multiplication*. It is important to note that while multiplication is well-defined mathematically, the process of multiplication calculation is not limited to traditional methods defined in textbook. Thus, examining only one calculation method presents a flawed experimental design that is vulnerable to exploitation. We selected four calculation methods that are representative to cover the major approaches to multiplication.

5.1 Historical and Modern Multiplication

In terms of multiplication, four different calculation methods are most representative from history to now: Standard Multiplication, Repetitive Addition, Lattice Method, and Egyptian Multiplication.

M1: Standard Multiplication In standard multiplication, we multiply each digit of one number by each digit of the other number, and then sum the results appropriately:

$$\begin{aligned} 12 \times 34 &= 12 \times (30 + 4) = 12 \times 30 + 12 \times 4 \\ &= 360 + 48 = 408 \end{aligned}$$

M2: Repetitive Addition Multiplication can be interpreted as repeated addition. For 12×34 , we add 12 thirty-four times:

$$\begin{aligned} 12 \times 34 &= 12 + 12 + 12 + \cdots + 12 \quad (34 \text{ times}) \\ &= 408 \end{aligned}$$

M3: Lattice Method In the lattice method (or grid method), we place the numbers along the top and side of a grid, perform single-digit multiplications, and then sum along the diagonals:

$$\begin{aligned} 12 \times 34 &= 10 \times 30 = 300 \\ &\quad 10 \times 4 = 40 \\ &\quad 2 \times 30 = 60 \\ &\quad 2 \times 4 = 8 \end{aligned}$$

Summing the results: $300 + 40 + 60 + 8 = 408$

M4: Egyptian Multiplication Egyptian multiplication computes the product by doubling the multiplicand and adding the results corresponding to the powers of two that sum to the multiplier. For 12×34 :

$$\begin{aligned} 12 \times 34 &= 12 \times 1 = 12 \\ 12 \times 2 &= \mathbf{24} \\ 12 \times 4 &= 48 \\ 12 \times 8 &= 96 \\ 12 \times 16 &= 192 \\ 12 \times 32 &= \mathbf{384} \end{aligned}$$

Summing the selected results: $24 + 384 = 408$

Since $34 = 2 + 32$, we select the results for 12×16 and 12×8 , and summing these gives the final product.

	Gemma-2-2B				Llama-3.1-8B			
	Standard	Lattice	Repetitive	Egyptian	Standard	Lattice	Repetitive	Egyptian
Task → Partial P.	+4.1%	+6.8%	-29.0%	+3.6%	+40.6%	+40.8%	-59.0%	+29.6%
Partial P. → Task	-6.1%	-10.7%	-20.3%	-9.6%	-3.7%	-0.2%	-0.9%	-2.7%

Table 1: Inductive and deductive accuracy difference Δ .

5.2 Examining Partial Product in Arithmetic Learning

To investigate whether LLMs generate partial products during arithmetic learning, we employ a set of diagnostic tasks as an approach to trace. We fine-tune Gemma-2-2B and Llama-3.1-8B on two-digit multiplication, observing changes in accuracy on diagnostic sets before and after fine-tuning (Task → Partial Products). Subsequently, we fine-tune the LLMs on these diagnostic sets and examine how their accuracy on the multiplication task changes (Partial Products → Task).

Method	Diagnostic Sets
Standard Multiplication	$\mathcal{P}_{\text{std}} = \{A_1 \times B_1 B_2, A_2 \times B_1 B_2, B_1 \times A_1 A_2, B_2 \times A_1 A_2\}$
Repetitive Addition	$\mathcal{P}_{\text{ra}} = \{\sum_{i=1}^{B_1 B_2} A_1 A_2, \sum_{i=1}^{A_1 A_2} B_1 B_2\}$
Lattice Method	$\mathcal{P}_{\text{lattice}} = \{A_1 0 \times B_1 0, A_1 0 \times B_2, A_2 \times B_1 0, A_2 \times B_2\}$
Egyptian Multiplication	$\mathcal{P}_{\text{egyptian}} = \{2^k \times A_1 A_2 \mid k \in 0, 1, \dots, [\log_2(B_1 B_2)]\}$

Table 2: Diagnostic sets with four calculation methods.

We probe language models’ partial product in four different directions. As shown in Table 2, for a task formatting like $A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4$, we would generate diagnostic test for each algorithm (detailed derivation in Appendix A.3).

Accuracy on Identifying Partial Products According to the results in Figure 2, we found that standard multiplication, the lattice method, and the Egyptian method significantly improved in identifying partial products after fine-tuning, with gains of +17.45%, +18.35%, and +10.45%, respectively. However, for repetitive addition tasks, LLMs failed to identify partial products, achieving an accuracy of only about 5% after fine-tuning.

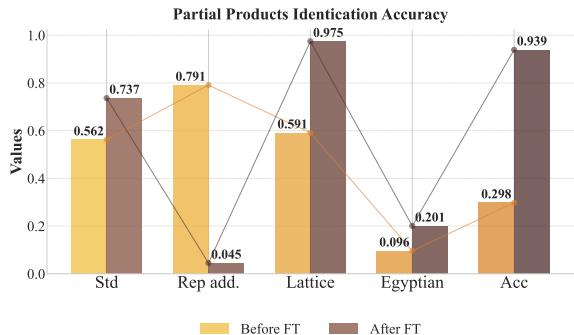


Figure 2: Partial products identification accuracy before and after fine-tuning on tasks. Scores are reported on average of Gemma-2-2B and Llama-3.1-8B.

A Deeper Look into Calculations Do the results showing increased accuracy across three paths really imply that partial products are used in arithmetic learning? We have two arguments against this interpretation. First, if LLMs genuinely use partial products to learn arithmetic, it is likely that they only use one calculation path at a time. Thus, the simultaneous improvement across three paths (standard, lattice, and Egyptian) is unusual. Second, if LLMs employ a specific path to compute partial products, this process should be demonstrated as bidirectional. Specifically, LLMs fine-tuned on a task should be able to identify partial products (inductive), and conversely, mastering partial products should enhance task learning (deductive). However, we currently have results for only one direction, lacking evidence for the other. Therefore, we extend our experiments to another direction.

Accuracy on Identifying Tasks We fine-tune two LLMs on diagnostic sets and present the results of identifying tasks before and after fine-tuning in Table 1. Our findings reveal that, fine-tuning specifically on partial products does not enhance task learning. Instead, it results in a performance drop across all four calculation paths for both models. This indicates that pre-learning partial products does not aid in arithmetic learning. The improved ability to recognize partial products appears to stem from the symbol learning process (note that the standard partial product $A_1 \times B_1 B_2$ is a sub-portion of $A_1 A_2 \times B_1 B_2$, similar to lattice and Egyptian methods) rather than being an intrinsic computational method used by the models.

6 Are Language Models Symbolic Observers?

An intuitive alternative for explaining increasing performance from inductive tasks (Task → Partial Products) is to treat LLMs as subgroup-level symbol observers, which aligns with the intrinsic properties of language modeling. Notably, the standard multiplication, lattice method, and Egyptian methods share a similar structure, where their partial product sets form token-level subgroups within the tasks. This observation naturally leads us to

explore this idea further.

6.1 Define Subgroup in Token Level

We first define subgroup in this section. Arithmetic learning involves a training set $\mathbb{D}_{\text{train}} = \{(a^{(k)}, b^{(k)}, c^{(k)})\}_{k=1}^N$ where $c^{(k)} = f(a^{(k)}, b^{(k)})$ and $f(\cdot)$ represents a binary operator, N is the number of dataset. In n -digit arithmetic learning, $a^{(k)}$ and $b^{(k)}$ can be regarded as different values taken by random variable sequences, $\{A_i\}_{i=1}^n$ and $\{B_i\}_{i=1}^n$, respectively. The random variables A_i and B_i all follow a discrete uniform distribution $P(X = x) = \frac{1}{10}, x = 0, 1, \dots, 9$. $c^{(k)}$ can be regarded as different values taken by random variable sequence $\{C_i\}_{i=1}^m$, where the random variables $C_i, i = 1, \dots, m$ has a joint distribution given by:

$$I_{\{f(a,b)=c\}} P(\{A_i\}_{i=1}^n = a) P(\{B_i\}_{i=1}^n = b)$$

where $I_{\{f(a,b)=c\}}$ is indicator function equals 1 if the condition $f(a, b) = c$ holds true, and 0 otherwise.

Definition 1 (Subgroup): For n -digit arithmetic, a subgroup s is defined as any element of the set \mathbb{S}_n :

$$s \in \mathbb{S}_n = \{((\mathbb{A}, \mathbb{B}), \mathbb{C}) \mid \mathbb{A} \subseteq \{A_i\}_{i=1}^n, \mathbb{B} \subseteq \{B_i\}_{i=1}^n, \mathbb{C} \subseteq \{C_i\}_{i=1}^m\}$$

where \mathbb{A} , \mathbb{B} and \mathbb{C} can be any subportion of random variable sequences $\{A_i\}_{i=1}^n$, $\{B_i\}_{i=1}^n$ and $\{C_i\}_{i=1}^m$, respectively. Specifically, we use $s_i \in \mathbb{S}_n$ to denote subgroups in i -th prediction for C_i .

Definition 2 (Subgroup Space): For any subgroup $s = ((\mathbb{A}, \mathbb{B}), \mathbb{C}) \in \mathbb{S}_n$, we have:

- Domain Space: $\mathcal{D}_s = \{\{\{a\}_{k=1}^{|\mathbb{A}|}, \{b\}_{k=1}^{|\mathbb{B}|}\} \mid P(\mathbb{A} = \{a\}_{k=1}^{|\mathbb{A}|}) > 0, P(\mathbb{B} = \{b\}_{k=1}^{|\mathbb{B}|}) > 0\}$. The size of domain space or cardinality is annotated as $|\mathcal{D}_s|$.
- Label Space: $\mathcal{L}_s = \{\{c\}_{k=1}^{|\mathbb{C}|} \mid P(\mathbb{C} = \{c\}_{k=1}^{|\mathbb{C}|}) > 0\}$. The size of label space is annotated as $|\mathcal{L}_s|$.

The entropy of label space is given by:

$$H(\mathcal{L}_s) = - \sum_{y \in \mathcal{L}_s} P(\mathbb{C} = y) \log_2 P(\mathbb{C} = y)$$

6.2 Difficulty in Arithmetic Learning: A Well-educated Hypothesis

We propose the following hypothesis: For an n -digit arithmetic task requiring m predictions, the overall difficulty ζ is related to:

$$\zeta \propto \hat{h}^m \quad (2)$$

where $\hat{h} = \prod_{i=1}^m h(s_i)^{\frac{1}{m}}$ represents the geometric average difficulty of an individual prediction, with $s_i \in \mathbb{S}_n$ denoting subgroup selection for the i -th prediction, and $h(\cdot)$ representing the subgroup complexity evaluation function.

Subgroup Complexity ($h(\cdot)$): This represents the most basic difficulty in arithmetic learning. It captures the inherent difficulty in the structure of the arithmetic learning tasks. We believe that the complexity on subgroup is strongly correlated with the property on that subgroup space:

- **Domain Space Cardinality $|\mathcal{D}|$:** The size of the domain space $|\mathcal{D}|$ determines how many data points are available for learning a pattern. If the label space is fixed, a larger domain space generally leads to improved learning outcomes.
- **Label Space Entropy $H(\mathcal{L})$:** Label space entropy $H(\mathcal{L})$ is also a critical factor in learning, as a low-entropy label space often leads to higher predictability.
- **Subgroup Quality Q :** For any subgroup $s = ((\mathbb{A}, \mathbb{B}), \mathbb{C}) \in \mathbb{S}_n$,

$$Q(s) = \max_{g \in \Omega_s} \sum_{a'=\{a\}_{k=1}^{|\mathbb{A}|}} \sum_{b'=\{b\}_{k=1}^{|\mathbb{B}|}} \sum_{c'=\{c\}_{k=1}^{|\mathbb{C}|}} P_s^f(a', b', c') P_s^p(g, a', b', c') \quad (3)$$

where $\Omega_s : \mathcal{D}_s \rightarrow \Theta(\mathcal{L}_s)$ represents a function space mapping from \mathcal{D}_s to $\Theta(\mathcal{L}_s)$, and $\Theta(\mathcal{L}_s)$ denotes the space of random variables taking values in \mathcal{L}_s . Here, $P_s^p(g, a', b', c') = P(g(a', b') = c')$ represents the probability that the function g maps (a', b') to c' , and $P_s^f(a', b', c')$ represents the probability that $\mathbb{A} = a'$, $\mathbb{B} = b'$ and $\mathbb{C} = c'$, hold simultaneously in arithmetic tasks. Thus, $Q(s)$ measures the maximum possible probability of predicting the value of \mathbb{C} that is consistent with the values in the dataset from the values of (\mathbb{A}, \mathbb{B}) .

Task	Format	C_1	C_2	C_3	C_4	C_5	$\{C_i\}_{i=1}^n$	
		$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$ \mathcal{L} $	$H(\mathcal{L})$
$f(a, b) = a + b$	$A_1 A_2 + B_1 B_2 = C_1 C_2 C_3$	0.9710	3.3215	3.3219	—	—	179	7.2130
$f(a, b) = a + b + 1$	$A_1 A_2 + B_1 B_2 = C_1 C_2 C_3$	0.9649	3.3215	3.3219	—	—	179	7.2130
$f(a, b) = a + b + 15$	$A_1 A_2 + B_1 B_2 = C_1 C_2 C_3$	0.9280	3.3214	3.3219	—	—	179	7.2130
$f(a, b) = a + b + 115$	$A_1 A_2 + B_1 B_2 = C_1 C_2 C_3$	0.9280	3.3214	3.3219	—	—	179	7.2130
$f(a, b) = (a + b) \text{ mod } 100$	$A_1 A_2 + B_1 B_2 = C_1 C_2$	3.3214	3.3219	—	—	—	100	6.6432
$f(a, b) = (a + b) \text{ mod } 50$	$A_1 A_2 + B_1 B_2 = C_1 C_2$	2.3217	3.3219	—	—	—	50	5.6436
$f(a, b) = (a + b) \text{ mod } 10$	$A_1 A_2 + B_1 B_2 = C_1$	3.3219	—	—	—	—	10	3.3219
$f(a, b) = a \times b$	$A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4$	2.8979	3.3215	3.3160	3.0340	—	2621	11.1172
$f(a, b) = a \times b \times 2$	$A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4 C_5$	0.6873	3.2173	3.3215	3.2964	2.2227	2621	11.1172
$f(a, b) = a \times b \times 4$	$A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4 C_5$	1.6030	3.3020	3.3204	3.2234	2.2227	2621	11.1172
$f(a, b) = a \times b \times 8$	$A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4 C_5$	2.5811	3.3202	3.3151	3.2235	2.2227	2621	11.1172
$f(a, b) = (a \times b) \text{ mod } 100$	$A_1 A_2 \times B_1 B_2 = C_1 C_2$	3.3160	3.0340	—	—	—	100	6.2912
$f(a, b) = (a \times b) \text{ mod } 50$	$A_1 A_2 \times B_1 B_2 = C_1 C_2$	2.3210	3.0340	—	—	—	50	5.3494
$f(a, b) = (a \times b) \text{ mod } 10$	$A_1 A_2 \times B_1 B_2 = C_1$	3.0340	—	—	—	—	10	3.0340

Table 3: Label space statistics with different rule perturbations. $H(\mathcal{L})$ represents the entropy of the label space, and $|\mathcal{L}|$ is the size of the label space. $\{C_j\}_{i=1}^n$ represents all positions in output digits.

Subgroup Selection (s_i): $s_i \in \mathbb{S}_n$ represents the subgroup selection for the i -th prediction. When LLMs predict the token in the i -th position, they must select subgroups that include C_i to align with the underlying pattern. This reflects the **learning dynamics** of language models in arithmetic tasks, abstractly linked to their decision-making and selection processes. As discussed in §6.4, LLMs seem to initially select the subgroup s_i with high quality $Q_{high}(s_i)$, progressing to lower quality $Q_{low}(s_i)$ (easy-to-hard) during learning.

6.3 Subgroup Complexity: Label Space Matters in the Final Stage

In this section, we discuss subgroup complexity in arithmetic learning. The domain space cardinality $|\mathcal{D}|$ represents the number of training data available, which is an obvious factor influencing learning. Subgroup quality $Q(s)$ will be detailed in §6.4. Thus, we primarily focus on *label space entropy* $H(\mathcal{L})$ in this section.

Rule Perturbation We first deliberately perturb the rules to observe whether these changes affect task difficulty for LLMs. We consider addition $f(a, b) = a + b$ and multiplication $f(a, b) = a \times b$ as our baselines. For addition, the perturbation is defined as $f(a, b) = a + b + \Delta c$, where $\Delta c = 1, 15, 115$ corresponds to perturbations at different position with different magnitudes. For multiplication, the perturbation is defined as $f(a, b) = a \times b \times \lambda$, where $\lambda = 2, 4, 8$ following similar reasons above. Additionally, we incorporate modular addition and multiplication as further perturbations. Table 3 showcases the basic change for label space entropy after applying perturbations. We then fine-tune Gemma-2-2B and Llama-3.1-8B using different perturbation rules to observe how well these large language models can be influenced from such perturbations in learning.

Results The results in Table 4 demonstrate that across two different rule perturbation methods and three distinct setups, both Gemma-2-2B and Llama-3.1-8B yield consistent outcomes. While calculating $13 \times 27 = 2808$ may seem counterintuitive, LLMs handle this arithmetic same as $13 \times 27 = 351$ when the label space entropy $H(\mathcal{L})$ are fixed.

	Gemma-2-2B	Llama-3.1-8B
$f(a, b) = a + b$	—	—
$f(a, b) = a + b + 1$	-0.1%	-0.1%
$f(a, b) = a + b + 15$	-0.9%	+0.1%
$f(a, b) = a + b + 115$	-1.4%	+0.7%
$f(a, b) = (a + b) \text{ mod } 100$	+10.1%	+3.7%
$f(a, b) = (a + b) \text{ mod } 50$	+13.1%	+6.7%
$f(a, b) = (a + b) \text{ mod } 10$	+26.1%	+13.7%
$f(a, b) = a \times b$	—	—
$f(a, b) = a \times b \times 2$	-1.1%	-2.7%
$f(a, b) = a \times b \times 4$	-1.7%	+0.7%
$f(a, b) = a \times b \times 8$	+0.2%	-3.7%
$f(a, b) = (a \times b) \text{ mod } 100$	+7.1%	+3.8%
$f(a, b) = (a \times b) \text{ mod } 50$	+12.1%	+5.3%
$f(a, b) = (a \times b) \text{ mod } 10$	+18.9%	+10.7%

Table 4: Test Accuracy difference Δ on perturbed addition and multiplication.

Regarding modular addition and multiplication with different modulus numbers, we find that decreasing the size of the entropy leads to performance improvements in both cases. These results highlight that an arithmetic task with low variability in the label space is more learnable. Together, these two observations reinforce the notion that LLMs are not performing traditional calculations but are instead functioning as sophisticated symbolic observers within the token space.

6.4 Subgroup Selection: Revealing Learning Dynamic in Arithmetic Learning

In the previous section, we established that subgroup space provides a basis for quantifying complexity in arithmetic tasks. However, the results mainly highlight insights at the end of learning

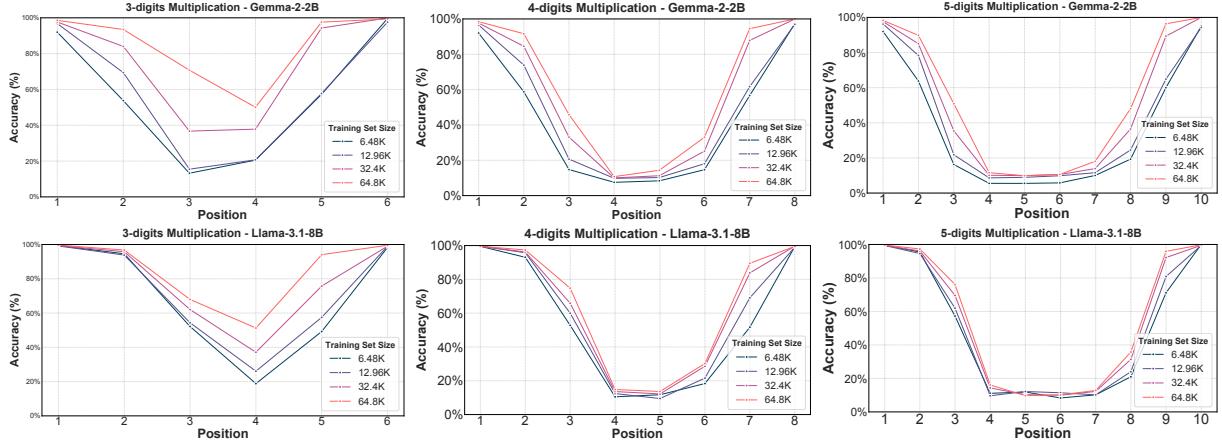


Figure 3: Position-level Accuracy from Gemma-2-2B and Llama-3.1-8B.

(test accuracy after 12 epochs), leaving the learning dynamics less explored. Here, we investigate these dynamics by analyzing digit-level accuracy in model outputs, observing how LLMs select subgroups S_n based on their performance across different positions.

Settings We maintain the same basic experimental settings as in the previous section to ensure the discussion remains within the same scope. We will train Gemma-2-2B and Llama-3.1-8B on four different dataset sizes ($6.48K$, $12.96K$, $32.4K$, and $64.8K$). Our experiments cover multiplication tasks ranging from 3-digit to 5-digit numbers, with output digits from 6 to 10.

Position-level Accuracy are U-curve Figure 3 reveals a phenomenon overlooked in previous studies. Contrary to the common assumption that position-level accuracy decreases from right to left due to carryover effects and least-to-most significant digit calculations (Lee et al., 2023; Zhang-Li et al., 2024), our results show a U-shaped accuracy curve in both Gemma-2-2B and Llama-3.1-8B models. Accuracy peaks at the beginning and end positions, exceeding 95%, with lower accuracy (~10%) in the middle positions, especially in higher-digit multiplication (e.g., 4th/5th for 4-digit, 5th/6th for 5-digit). These results provide valuable insights, suggesting that the difficulty in learning multiplication is concentrated in the middle positions rather than at the beginning, which conceptually corresponds to the final steps of calculation.

Subgroup Selection via Quality We think that the U-curve occurs because the subgroup for middle digits has a lower subgroup quality $Q(s)$ compared to the beginning and end digits. Given a subgroup $s = ((\mathbb{A}, \mathbb{B}), \mathbb{C})$, where \mathbb{C}_{mid} represents

the middle digits, more digits from the operands are required to determine the value of \mathbb{C}_{mid} compared to representing digits at the beginning or end positions. This leads to a larger domain size $|\mathcal{D}|$ when summing all candidates. Consequently, determining the value of \mathbb{C} becomes less certain, resulting in a lower $Q(s)$. Specifically, when \mathbb{C} represents the last digit, its value can be fully determined by the least significant digits of the operands, which is not the case when \mathbb{C} represents a middle digit. For instance, in the case of 3-digit multiplication, it is relatively easy to identify the subgroup $s = \{(\{\mathbb{A}_1\}, \{\mathbb{B}_1\}), \{\mathbb{C}_1\}\}$ for learning the first digits or the subgroup $s = \{(\{\mathbb{A}_3\}, \{\mathbb{B}_3\}), \{\mathbb{C}_6\}\}$ for learning the last digits. This observation explains why the digits at the beginning or the end are easier to learn. It also reveals that LLMs fit arithmetic learning patterns following an easy-to-hard mechanism (from high $Q(s)$ to low $Q(s)$), which demonstrates gradient descent in the "fastest" symbolic direction.

7 Conclusions

In this work, we investigate whether LLMs solve arithmetic tasks using partial products or operate as symbolic pattern matchers. Our findings confirm that LLMs do not rely on partial products but instead approach arithmetic purely symbolically. By breaking tasks into subgroups, we demonstrate that the difficulty in arithmetic learning can be attributed to subgroup complexity and selection. Our results emphasize the crucial role of label space entropy in understanding the convergence stage and the quality of subgroups for learning dynamics. Overall, at least in our setting, LLMs function as purely symbolic learners in arithmetic tasks. We encourage further research to explore more complex tasks from this perspective.

8 Limitations

In terms of limitations, one area our work does not currently address is the application of our framework to different chain-of-thought (CoT) methods (Wei et al., 2023; Deng et al., 2024). While CoT has proven to be highly effective in arithmetic learning, particularly by decomposing overall difficulty into smaller, more manageable operations—thereby reducing subgroup complexity from exponential to linear—this aspect has not been explored in our study. Additionally, we have not applied our framework in a totally natural language-aware setting like GSM8K or MATH. Exploring how LLMs leverage their symbolic capabilities in such a context could provide deeper insights into their reasoning abilities, particularly in tasks that require structured, multi-step problem solving. These unexplored areas present significant opportunities for future research.

9 Ethics Statement

Our research primarily focuses on the symbolic learning capabilities of large language models in arithmetic tasks. As such, it does not involve the collection or use of any human data. No personal or sensitive information is handled or analyzed in this study. We acknowledge the potential biases inherent in the datasets used for model training and the limitations of relying on symbolic learning without fully understanding the underlying numerical or logical processes. The societal impact of increased reliance on LLMs for arithmetic tasks, including overconfidence in symbolic learning without full comprehension, warrants careful consideration. We advocate for transparent model evaluations and awareness of the limitations in deploying such models for critical decision-making.

References

- Anthropic. 2024. *Claude*.
- Daman Arora, Himanshu Gaurav Singh, and Mausam. 2023. [Have llms advanced enough? a challenging problem solving benchmark for large language models](#). *ArXiv*, abs/2305.15074.
- Sasha Boguraev, Ben Lipkin, Leonie Weissweiler, and Kyle Mahowald. 2024. [Models can and should embrace the communicative nature of human-generated math](#). *Preprint*, arXiv:2409.17005.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *Preprint*, arXiv:2110.14168.
- Yuntian Deng, Yejin Choi, and Stuart Shieber. 2024. [From explicit cot to implicit cot: Learning to internalize cot step by step](#). *Preprint*, arXiv:2405.14838.
- Yuntian Deng, Kiran Prasad, Roland Fernandez, Paul Smolensky, Vishrav Chaudhary, and Stuart Shieber. 2023. [Implicit chain of thought reasoning via knowledge distillation](#). *Preprint*, arXiv:2311.01460.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, and et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Siyuan Guo, Aniket Didolkar, Nan Rosemary Ke, Anirudh Goyal, Ferenc Huszár, and Bernhard Schölkopf. 2024. [Learning beyond pattern matching? assaying mathematical understanding in llms](#). *Preprint*, arXiv:2405.15485.
- Michael Hanna, Ollie Liu, and Alexandre Variengien. 2023. [How does gpt-2 compute greater-than?: Interpreting mathematical abilities in a pre-trained language model](#). *Preprint*, arXiv:2305.00586.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. [Measuring mathematical problem solving with the math dataset](#). *Preprint*, arXiv:2103.03874.
- Danqing Huang, Shuming Shi, Chin-Yew Lin, Jian Yin, and Wei-Ying Ma. 2016. [How well do computers solve math word problems? large-scale dataset construction and evaluation](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 887–896, Berlin, Germany. Association for Computational Linguistics.
- Damjan Kalajdzievski. 2023. [A rank stabilization scaling factor for fine-tuning with lora](#). *Preprint*, arXiv:2312.03732.
- Andrej Karpathy. 2022. Andrej karpathy’s lightweight implementation of medium-sized gpts. <https://github.com/karpathy/nanoGPT>. Accessed: 2024-09-28.
- Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. 2014. [Learning to automatically solve algebra word problems](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages

- 271–281, Baltimore, Maryland. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Nayoung Lee, Kartik Sreenivasan, Jason D. Lee, Kangwook Lee, and Dimitris Papailiopoulos. 2023. [Teaching arithmetic to small transformers](#). *Preprint*, arXiv:2307.03381.
- Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. [Solving quantitative reasoning problems with language models](#). In *Advances in Neural Information Processing Systems*.
- Yixin Liu, Avi Singh, C. Daniel Freeman, John D. Co-Reyes, and Peter J. Liu. 2023. [Improving large language model fine-tuning for solving math problems](#). *Preprint*, arXiv:2310.10047.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, and et al. 2024. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Alessandro Stolfo, Yonatan Belinkov, and Mrinmaya Sachan. 2023. [A mechanistic interpretation of arithmetic reasoning in language models using causal mediation analysis](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7035–7052, Singapore. Association for Computational Linguistics.
- Sowmya S. Sundaram, Sairam Gurajada, Marco Fisichella, Deepak P, and Savitha Sam Abraham. 2022. [Why are nlp models fumbling at elementary math? a survey of deep learning based word problem solvers](#). *ArXiv*, abs/2205.15683.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, and et al. 2024. [Gemma 2: Improving open language models at a practical size](#). *Preprint*, arXiv:2408.00118.
- Avijit Thawani, Jay Pujara, Filip Ilievski, and Pedro Szekely. 2021. [Representing numbers in NLP: a survey and a vision](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–656, Online. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). *Preprint*, arXiv:2203.11171.
- Yan Wang, Xiaojiang Liu, and Shuming Shi. 2017. [Deep neural solver for math word problems](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 845–854, Copenhagen, Denmark. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. [Chain-of-thought prompting elicits reasoning in large language models](#). *Preprint*, arXiv:2201.11903.
- Yiran Wu, Feiran Jia, Shaokun Zhang, Hangyu Li, Erkang Zhu, Yue Wang, Yin Tat Lee, Richard Peng, Qingyun Wu, and Chi Wang. 2024a. [Mathchat: Converse to tackle challenging math problems with llm agents](#). *Preprint*, arXiv:2306.01337.
- Zhengxuan Wu, Atticus Geiger, Thomas Icard, Christopher Potts, and Noah D. Goodman. 2024b. [Interpretability at scale: Identifying causal mechanisms in alpaca](#). *Preprint*, arXiv:2305.08809.
- Roy Xie, Chengxuan Huang, Junlin Wang, and Bhuvan Dhingra. 2024. [Adversarial math word problem generation](#). *Preprint*, arXiv:2402.17916.
- Zhen Yang, Ming Ding, Qingsong Lv, Zhihuan Jiang, Zehai He, Yuyi Guo, Jinfeng Bai, and Jie Tang. 2023. [Gpt can solve mathematical problems without a calculator](#). *Preprint*, arXiv:2309.03241.
- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2023. [Mammoth: Building math generalist models through hybrid instruction tuning](#). *ArXiv*, abs/2309.05653.
- Wei Zhang, Chaoqun Wan, Yonggang Zhang, Yiu ming Cheung, Xinmei Tian, Xu Shen, and Jieping Ye. 2024. [Interpreting and improving large language models in arithmetic calculation](#). *Preprint*, arXiv:2409.01659.
- Daniel Zhang-Li, Nianyi Lin, Jifan Yu, Zheyuan Zhang, Zijun Yao, Xiaokang Zhang, Lei Hou, Jing Zhang,

and Juanzi Li. 2024. Reverse that number! decoding order matters in arithmetic learning. *Preprint*, arXiv:2403.05845.

Yilun Zhao, Hongjun Liu, Yitao Long, Rui Zhang, Chen Zhao, and Arman Cohan. 2024a. Financemath: Knowledge-intensive math reasoning in finance domains. *Preprint*, arXiv:2311.09797.

Yilun Zhao, Yitao Long, Hongjun Liu, Ryo Kamoi, Linyong Nan, Lyuhao Chen, Yixin Liu, Xiangru Tang, Rui Zhang, and Arman Cohan. 2024b. Docmath-eval: Evaluating math reasoning capabilities of llms in understanding long and specialized documents. *Preprint*, arXiv:2311.09805.

Andy Zhou, Kai Yan, Michal Shlapentokh-Rothman, Haohan Wang, and Yu-Xiong Wang. 2023a. Language agent tree search unifies reasoning acting and planning in language models. *ArXiv*, abs/2310.04406.

Zihao Zhou, Qifeng Wang, Mingyu Jin, Jie Yao, Jianan Ye, Wei Liu, Wei Wang, Xiaowei Huang, and Kaizhu Huang. 2023b. Mathattack: Attacking large language models towards math solving ability. *Preprint*, arXiv:2309.01686.

A Appendix

A.1 Training Detail

We carefully tuned the hyperparameters in our experiments. The learning rate for Gemma-2-2B was set to $1e - 4$, while for Llama-3.1-8B it was $2e - 4$. Both models used a warm-up of 5 steps and a weight decay of 0.01. We trained for 12 epochs, splitting the dataset into 80% for training, 10% for validation, and 10% for testing. Evaluation was conducted at the end of each epoch, with checkpoints saved based on the best performance on the validation set. LoRA finetuning is used for both models with same setting, lora_rank = 64, lora_alpha = 16, lora_dropout = 0, and we enable rank stabilized lora (Kalajdzievski, 2023) during training. We use unsloth¹ to fasten our training process and vLLM (Kwon et al., 2023) to increase our inference process.

A.2 Format Perturbations in Arithmetic Tasks

In this section, we apply three types of format perturbations to basic addition and multiplication tasks to evaluate the symbolic reasoning capabilities of large language models (LLMs). Our experiments utilize Gemma-2-2B and Llama-3.1-8B models. The primary objective of varying the input format

is to investigate whether LLMs function as purely symbolic learners in arithmetic tasks. We consider the following three types of format perturbation in Table 5:

1. **Natural Language (NL):** Arithmetic expressions are converted into natural language statements. For example, the equation “3 + 5” becomes “What is 3 times 5?”
2. **Random String (RS):** Arithmetic expressions are first converted into natural language and then replaced with meaningless strings. Using the previous example, “3 + 5” might be transformed into “flad kf 3 lfd 5?”
3. **Disturbed Digits (DD):** Arithmetic expressions are initially converted into natural language and subsequently replaced with random digits. For instance, “3 + 5” could become “65.1 44 3 4 5?” This approach creates a counterfactual context for arithmetic tasks, increasing the difficulty for the models.

By implementing these perturbations, we aim to assess the robustness of LLMs in handling arithmetic operations under varied and challenging input formats.

Results We examined the impact of three types of input format perturbations (Natural Language (NL), Random String (RS), and Disturbed Digits (DD)) on the arithmetic reasoning tasks for Gemma-2.2B and Llama-3.1-8B models. Table 6 shows that across both addition and multiplication tasks, the performance of the models remains largely unaffected by the perturbations when the label space is fixed. Specifically, there is only a marginal change in accuracy under the NL and RS formats, while the DD format causes minor fluctuations but does not significantly degrade performance. This demonstrates that LLMs can effectively handle various input perturbations as long as the output space remains consistent, suggesting their robustness in symbolic reasoning tasks despite superficial input variations.

A.3 Mathematical Explanation of Diagnostic Sets for Multiplication Algorithms

For the multiplication task involving two two-digit numbers formatted as $A_1A_2 \times B_1B_2 = C_1C_2C_3C_4$, we generate diagnostic test sets \mathcal{P} for each algorithm to analyze and understand the partial computations involved. Below, we provide a

¹Unsloth AI is at: <https://unsloth.ai/>

Task	Format	C_1	C_2	C_3	C_4	C_5	$\{C_j\}_{i=1}^n$	
		$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$H(\mathcal{L})$	$ \mathcal{L} $	$H(\mathcal{L})$
$f(a, b) = a + b$	$A_1 A_2 + B_1 B_2 = C_1 C_2 C_3$	0.9710	3.3215	3.3219	—	—	179	7.2130
$f(a, b) = a + b$	What is $A_1 A_2$ add $B_1 B_2$? Answer: $C_1 C_2 C_3$	0.9649	3.3215	3.3219	—	—	179	7.2130
$f(a, b) = a + b$	fafr if $A_1 A_2$ hfk $B_1 B_2$? Ffhjar: $C_1 C_2 C_3$	3.3214	3.3219	—	—	—	179	7.2130
$f(a, b) = a + b$	3.123 34 $A_1 A_2$ 461 $B_1 B_2$? 952414: $C_1 C_2 C_3$	0.9280	3.3214	3.3219	—	—	179	7.2130
$f(a, b) = a \times b$	$A_1 A_2 \times B_1 B_2 = C_1 C_2 C_3 C_4 C_5$	2.5811	3.3202	3.3151	3.2235	2.2227	2621	11.1172
$f(a, b) = a \times b$	What is $A_1 A_2$ multiply $B_1 B_2$? Answer: $C_1 C_2 C_3 C_4$	2.8979	3.3215	3.3160	3.0340	—	2621	11.1172
$f(a, b) = a \times b$	fafr if $A_1 A_2$ hfk $B_1 B_2$? Ffhjar: $C_1 C_2 C_3 C_4$	0.6873	3.2173	3.3215	3.2964	2.2227	2621	11.1172
$f(a, b) = a \times b$	3.123 34 $A_1 A_2$ 461 $B_1 B_2$? 952414: $C_1 C_2 C_3 C_4$	1.6030	3.3020	3.3204	3.2234	2.2227	2621	11.1172

Table 5: Label space statistics with different format perturbations. $H(\mathcal{L})$ represents the entropy of the space, and $|\mathcal{L}|$ is the size of the space. $\{C_j\}_{i=1}^n$ represents all possible output digits.

	Format	Gemma-2-2B	Llama-3.1-8B
$f(a, b) = a + b$	Natural Language	—	—
$f(a, b) = a + b$	Random String	+0.1%	-0.2%
$f(a, b) = a + b$	Disturbed Digits	-3.9%	-2.1%
$f(a, b) = a \times b$	Natural Language	—	—
$f(a, b) = a \times b$	Random String	+0.3%	-0.5%
$f(a, b) = a \times b$	Disturbed Digits	-1.9%	-3.1%

Table 6: Test Accuracy difference Δ on perturbed addition and multiplication.

mathematical explanation for the formulation of these diagnostic sets for each multiplication algorithm.

A.3.1 Standard Multiplication

In the standard multiplication algorithm, we multiply each digit of one number by each digit of the other number and sum the appropriately weighted results.

Formulation: Let the two-digit numbers be expressed as:

$$\begin{aligned} a &= A_1 A_2 = 10A_1 + A_2, \\ b &= B_1 B_2 = 10B_1 + B_2. \end{aligned}$$

The product is:

$$ab = (10A_1 + A_2)(10B_1 + B_2).$$

Expanding, we get four partial products:

$$ab = 100A_1 B_1 + 10A_1 B_2 + 10A_2 B_1 + A_2 B_2.$$

Diagnostic Set:

$$\mathcal{P}_{\text{std}} = \{A_1 \times b, A_2 \times b, B_1 \times a, B_2 \times a\}.$$

Explanation:

- $A_1 \times b$: Multiplying the tens digit of a by the entire number b :

$$A_1 \times b = A_1 \times (10B_1 + B_2) = 10A_1 B_1 + A_1 B_2.$$

- $A_2 \times b$: Multiplying the units digit of a by b :

$$A_2 \times b = A_2 \times (10B_1 + B_2) = 10A_2 B_1 + A_2 B_2.$$

- $B_1 \times a$: Multiplying the tens digit of b by a :

$$B_1 \times a = B_1 \times (10A_1 + A_2) = 10A_1 B_1 + A_2 B_1.$$

- $B_2 \times a$: Multiplying the units digit of b by a :

$$B_2 \times a = B_2 \times (10A_1 + A_2) = 10A_1 B_2 + A_2 B_2.$$

Including these partial products in \mathcal{P}_{std} captures all intermediary computations in the standard algorithm, facilitating a comprehensive diagnostic analysis.

A.3.2 Repetitive Addition

Repetitive addition interprets multiplication as adding one number to itself repeatedly.

Diagnostic Set:

$$\mathcal{P}_{\text{ra}} = \left\{ \sum_{i=1}^b a, \sum_{j=1}^a b \right\}.$$

Explanation:

- $\sum_{i=1}^b a$: Adding a to itself b times:

$$\sum_{i=1}^b a = a + a + \cdots + a \quad (b \text{ times}) = ab.$$

- $\sum_{j=1}^a b$: Adding b to itself a times:

$$\sum_{j=1}^a b = b + b + \cdots + b \quad (a \text{ times}) = ab.$$

Both summations lead to the same product ab , and including them in \mathcal{P}_{ra} allows for analyzing both repetitive addition paths in the algorithm.

A.3.3 Lattice Method

The lattice method (or grid method) organizes the multiplication of each digit pair in a grid and sums along diagonals.

Diagnostic Set:

$$\mathcal{P}_{\text{lattice}} = \{A_1 \times B_1, A_1 \times B_2, A_2 \times B_1, A_2 \times B_2\}.$$

Explanation:

- $A_1 \times B_1$: Tens digit of a times tens digit of b .
- $A_1 \times B_2$: Tens digit of a times units digit of b .
- $A_2 \times B_1$: Units digit of a times tens digit of b .
- $A_2 \times B_2$: Units digit of a times units digit of b .

These products fill the cells of the lattice grid:

	B_1	B_2
A_1	A_1B_1	A_1B_2
A_2	A_2B_1	A_2B_2

Summing along the diagonals yields the final product. Including these partial products in $\mathcal{P}_{\text{lattice}}$ covers all the necessary computations in the lattice method.

A.3.4 Egyptian Multiplication

Egyptian multiplication involves doubling the multiplicand and adding specific results based on the binary representation of the multiplier.

Diagnostic Set:

$$\mathcal{P}_{\text{egyptian}} = \left\{ 2^k \times a \mid k = 0, 1, \dots, \lfloor \log_2 b \rfloor \right\}.$$

Explanation:

- **Binary Representation of b :** Express b as a sum of powers of two:

$$b = \sum_{k=0}^n C_k 2^k, \quad C_k \in \{0, 1\}, \quad n = \lfloor \log_2 b \rfloor.$$

- **Doubling a :** Compute successive doublings of a :

$$2^0 \times a, 2^1 \times a, \dots, 2^n \times a.$$

- **Selection and Summation:** Identify which $2^k \times a$ correspond to $C_k = 1$ in b 's binary representation and sum them:

$$ab = \sum_{k=0}^n C_k (2^k \times a).$$

Including all $2^k \times a$ up to n in $\mathcal{P}_{\text{egyptian}}$ ensures that we have the necessary partial products for any b , allowing us to reconstruct ab by selecting and summing the appropriate terms.

Example:

If $b = 13$, its binary representation is 1101, so $b = 2^3 + 2^2 + 2^0$. The partial products are:

$$\begin{aligned} 2^0 \times a &= a, \\ 2^1 \times a &= 2a, \\ 2^2 \times a &= 4a, \\ 2^3 \times a &= 8a. \end{aligned}$$

Select $2^0 \times a$, $2^2 \times a$, and $2^3 \times a$ (since $C_0 = C_2 = C_3 = 1$) and sum:

$$ab = a + 4a + 8a = 13a.$$

By formulating the diagnostic sets \mathcal{P} as above for each multiplication algorithm, we encapsulate all intermediary computational steps inherent to each method.